**STUDENT**

Xin Qu

**AUTOGRADER SCORE**

**81.92 / 100.0**

## held_out_accu vs Steps



Legend:
- $\lambda=1e-3$
- $\lambda=1e-2$
- $\lambda=1e-1$
- $\lambda=1$

## Validation accuracy vs Steps



Legend:
- $\lambda=1e-3$
- $\lambda=1e-2$
- $\lambda=1e-1$
- $\lambda=1$

magnitude vs Steps

Based on the accuracy score of validation set (10% of training set), the highest accuracy score is 0.80 when lambda = 1e-3.

```
1   accuracy = []
2   magnitude = []
3   a_list = []
4   b_list = []
5   test_accu_list = []
6   held_out = []
7   for i in lambdas:
8       cur_acc, cur_mag, a, b, test_accu, held_out_accu = SVM(train, i)
9       accuracy.append(cur_acc)
10      magnitude.append(cur_mag)
11      a_list.append(a)
12      b_list.append(b)
13      test_accu_list.append(test_accu)
14      held_out.append(held_out_accu)
```

```
1   max(test_accu_list)
```

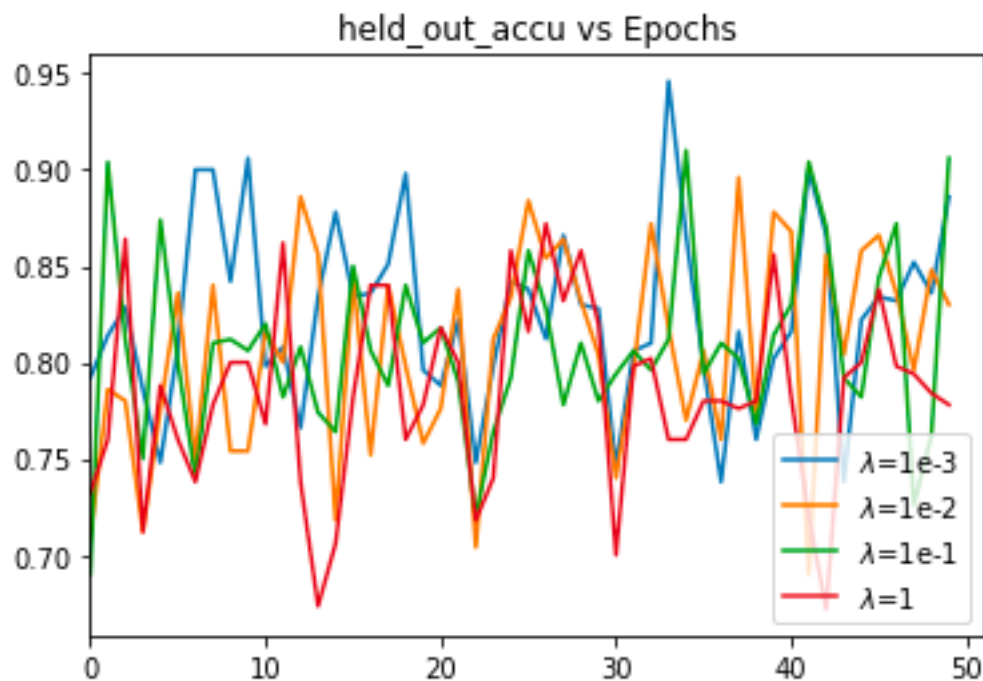0.8002729754322111

```
1   test_accu_list = np.array(test_accu_list)
2   np.amax(test_accu_list), np.argmax(test_accu_list)
```

(0.8002729754322111, 0)

By learning rate, I set up steplength (learning rate) = 1 / (0.01 * epoch + 50). From the following figure, when epoch reaches around 32 and lambda = 1e-3, the accuracy score for held out epoch reaches its highest value.

page 5:

1. SVM training with stochastic gradient descent updating

```
 1  def SVM(traindata, lamb):
 2      train_data, test_data = train_test_split(train)
 3      accuracy_list = []
 4      magnitude_list = []
 5      a = np.ones((1, 6))
 6      b = 1
 7      held_out_accu = []
 8      ##50 epochs at least 300 steps
 9      for epoch in range(50): ### 50 epochs
10          steplength = 1 / (0.01 * epoch + 50)
11          ###random seperate 50 examples
12          held_out_index = set(np.random.choice(train_data.shape[0], size = 50))
13          epoch_index = set(range(train_data.shape[0])) - held_out_index
14          held_out_index = list(held_out_index)
15          held_out = train_data[held_out_index]
16          epoch_data = train_data[list(epoch_index)]
17          batch_size = 1
18          #held_out_accu = []
19          for s in range(300):
20              ##ramdon choose batch
21              batch_index = np.random.choice(held_out.shape[0], batch_size)
22              batch_data = held_out[batch_index]
23              boundary = np.dot(batch_data[:, -1], (np.dot(batch_data[:, 0:-1], a.T) + b))
24              if boundary.item(0) >= 1:
25                  a = a - steplength * lamb * a
26                  b = b
27              else:
28                  a = a - steplength * (lamb * a - batch_data[:, 0:-1] * batch_data[:, -1])
29                  b = b - steplength * (-batch_data[:, -1])
30              if s % 30 == 0:
31                  ###accuracy score of validation set and magnitude every 30 steps
32                  accuracy_list.append(accuracyscore(predict(test_data, a, b), np.ravel(test_data[:, -1])))
33                  magnitude_list.append(np.dot(a, a.T).item(0))
34                  held_out_accu.append(accuracyscore(predict(held_out, a, b), np.ravel(held_out[:, -1])))
35      test_accu = accuracyscore(predict(test_data, a, b), np.ravel(test_data[:, -1]))
36      return accuracy_list, magnitude_list, a, b, test_accu, held_out_accu
```

2. Label prediction

```
def predict(x, a, b):###predict train data with class label
    result = np.dot(x[:, 0:-1], a.T) + b
    result[result >=0] = 1
    result[result < 0] = -1
    return np.ravel(result)
```

3. Calculation of the accuracies

```
def accuracyscore(x, y):
    return sum(x == np.ravel(y)) / float(len(y))
```

# Page 6 All Codes

```python
In [1]:   1  import pandas as pd
          2  import numpy as np
          3  import matplotlib.pyplot as plt
          4  train = pd.read_csv('/Users/xinqu/Sandbox/CS498 Applied Machine Learning/HW/HW2/train.txt', header = None).values
          5  test = pd.read_csv('/Users/xinqu/Sandbox/CS498 Applied Machine Learning/HW/HW2/test.txt', header = None).values
```

```python
In [2]:   1  ###label column14
          2  train[np.where(train == ' >50K')] = 1
          3  train[np.where(train == ' <=50K')] = -1
          4  train = train[:, [0, 2, 4, 10, 11, 12, 14]] ###keep continuous attributes
          5  ##scale train data
          6  from sklearn.preprocessing import scale
          7  train[:, 0:-1] = scale(train[:, 0:-1].astype(float), with_mean = True) ###0 mean, 1 std
```

```python
In [3]:   1  lambdas = np.array([0.001, 0.01, 0.1, 1])
```

```python
In [4]:   1  import random
          2  def train_test_split(data, ratio = 0.1):
          3      ###return train_data, test_data, split test_data by default ratio = 0.1
          4      random.shuffle(data)
          5      train_size = int((1 - ratio) * data.shape[0])
          6      train_data = data[:train_size]
          7      test_data = data[train_size:]
          8      return train_data, test_data
          9  def accuracyscore(x, y):
         10      return sum(x == np.ravel(y)) / float(len(y))
         11  #train_test_slit(train)
         12  def predict(x, a, b):###predict train data with class label
         13      result = np.dot(x[:, 0:-1], a.T) + b
         14      result[result >=0] = 1
         15      result[result < 0] = -1
         16      return np.ravel(result)
```

```python
In [5]:   1  def SVM(traindata, lamb):
          2      train_data, test_data = train_test_split(train)
          3      accuracy_list = []
          4      magnitude_list = []
          5      a = np.ones((1, 6))
          6      b = 1
          7      held_out_accu = []
          8      ##50 epochs at least 300 steps
          9      for epoch in range(50): ### 50 epochs
         10          steplength = 1 / (0.01 * epoch + 50)
         11          ###random seperate 50 examples
         12          held_out_index = set(np.random.choice(train_data.shape[0], size = 50))
         13          epoch_index = set(range(train_data.shape[0])) - held_out_index
         14          held_out_index = list(held_out_index)
         15          held_out = train_data[held_out_index]
         16          epoch_data = train_data[list(epoch_index)]
         17          batch_size = 1
         18          #held_out_accu = []
         19          for s in range(300):
         20              ##ramdon choose batch
         21              batch_index = np.random.choice(held_out.shape[0], batch_size)
         22              batch_data = held_out[batch_index]
         23              boundary = np.dot(batch_data[:, -1], (np.dot(batch_data[:, 0:-1], a.T) + b))
         24              if boundary.item(0) >= 1:
         25                  a = a - steplength * lamb * a
         26                  b = b
         27              else:
         28                  a = a - steplength * (lamb * a - batch_data[:, 0:-1] * batch_data[:, -1])
         29                  b = b - steplength * (-batch_data[:, -1])
         30              if s % 30 == 0:
         31                  ###accuracy score of validation set and magnitude every 30 steps
         32                  accuracy_list.append(accuracyscore(predict(test_data, a, b), np.ravel(test_data[:, -1])))
         33                  magnitude_list.append(np.dot(a, a.T).item(0))
         34                  held_out_accu.append(accuracyscore(predict(held_out, a, b), np.ravel(held_out[:, -1])))
         35      test_accu = accuracyscore(predict(test_data, a, b), np.ravel(test_data[:, -1]))
         36      return accuracy_list, magnitude_list, a, b, test_accu, held_out_accu
```

```
In [6]:   1  accuracy = []
          2  magnitude = []
          3  a_list = []
          4  b_list = []
          5  test_accu_list = []
          6  held_out = []
          7  for i in lambdas:
          8      cur_acc, cur_mag, a, b, test_accu, held_out_accu = SVM(train, i)
          9      accuracy.append(cur_acc)
         10      magnitude.append(cur_mag)
         11      a_list.append(a)
         12      b_list.append(b)
         13      test_accu_list.append(test_accu)
         14      held_out.append(held_out_accu)
```
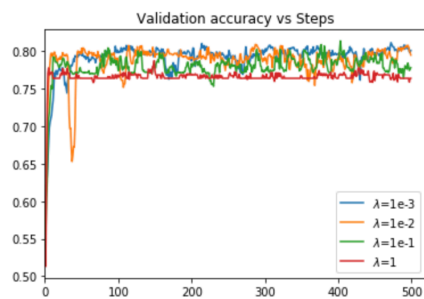
```
In [7]:   1  max(test_accu_list)
```
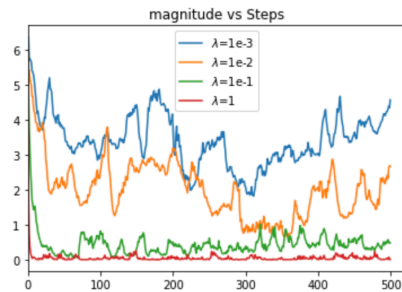
Out[7]:  0.8002729754322111

```
In [8]:   1  test_accu_list = np.array(test_accu_list)
          2  np.amax(test_accu_list), np.argmax(test_accu_list)
```

Out[8]:  (0.8002729754322111, 0)

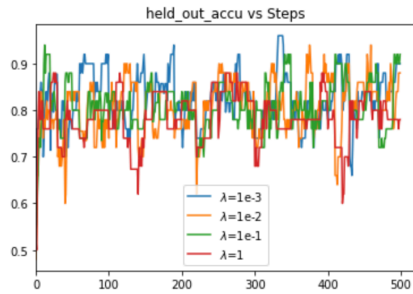```
In [9]:   1  steps = np.arange(500)
          2  plt.figure(1)
          3  for i in range(4):
          4      plt.plot(steps, accuracy[i])
          5  plt.xlim(0, 520)
          6  plt.legend(['$\lambda$=1e-3', '$\lambda$=1e-2', '$\lambda$=1e-1', '$\lambda$=1'], loc='lower right')
          7  plt.title('Validation accuracy vs Steps')
          8  plt.show()
```



Validation accuracy vs Steps
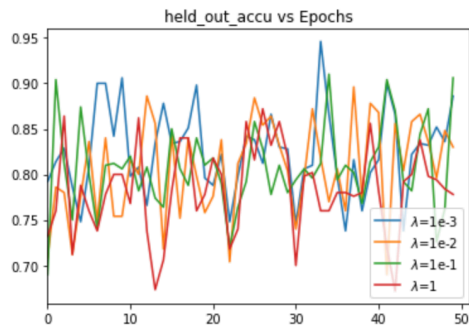
```
In [10]:  1  plt.figure(2)
          2  for i in range(4):
          3      plt.plot(steps, magnitude[i])
          4  plt.xlim(0, 520)
          5  plt.legend(['$\lambda$=1e-3', '$\lambda$=1e-2', '$\lambda$=1e-1', '$\lambda$=1'], loc='best')
          6  plt.title('magnitude vs Steps')
          7  plt.show()
```



magnitude vs Steps

```
In [12]:  1  plt.figure(3)
          2  for i in range(4):
          3      plt.plot(steps, held_out[i])
          4  plt.xlim(0, 520)
          5  plt.legend(['$\lambda$=1e-3', '$\lambda$=1e-2', '$\lambda$=1e-1', '$\lambda$=1'], loc='best')
          6  plt.title('held_out_accu vs Steps')
          7  plt.show()
```


held_out_accu vs Steps

```
In [18]:   1  plt.figure(4)
           2  epochstep = np.arange(50)
           3  for i in range(4):
           4      accuracy[i] = np.array(held_out[i])
           5      tmp = accuracy[i].reshape((50, 10))
           6      avg = np.mean(tmp, axis = 1)
           7      plt.plot(epochstep, avg)
           8  plt.xlim(0, 51)
           9  plt.legend(['$\lambda$=1e-3', '$\lambda$=1e-2', '$\lambda$=1e-1', '$\lambda$=1'], loc='best')
          10  plt.title('held_out_accu vs Epochs')
          11  plt.show()
```


held_out_accu vs Epochs

## The best value of regularization constant is 1e-3

```
In [13]:  1  test = test[:, [0, 2, 4, 10, 11, 12]] ###keep continuous attributes
          2  ##scale test data
          3  test[:,:] = scale(test.astype(float), with_mean = True) ###0 mean, 1 std
```

```
In [14]:  1  def predicty(y, a, b):###predict data without class label
          2      result = np.dot(y, a.T) + b
          3      result[result >=0] = 1
          4      result[result < 0] = -1
          5      return np.ravel(result)
          6  test_y_pred = predicty(test, a_list[0], b_list[0])
```

```
In [15]:  1  from collections import Counter
          2  Counter(test_y_pred)
```

Out[15]:  Counter({-1: 4011, 1: 874})

```
In [16]:  1  test_y_class = []
          2  for i in range(len(test_y_pred)):
          3      if test_y_pred[i] == 1:
          4          test_y_class.append('>50K')
          5      else:
          6          test_y_class.append('<=50K')
          7  with open('submission.txt', 'w') as f:
          8      for item in test_y_class:
          9          f.write("%s\n" % item)
```